

# INFRASTRUCTURE OPTIMIZATION IN KUBERNETES CLUSTER

Darshan Haragi L<sup>1</sup>, Mahith S<sup>2</sup>, Prof. Sahana B<sup>3</sup>

<sup>1</sup>Student, Dept of Electronics and Communication, R.V. College of Engineering, Bangalore, INDIA -560059

<sup>2</sup>Student, Dept of Electronics and Communication, R.V. College of Engineering, Bangalore, INDIA -560059

<sup>3</sup>Assistant Professor, Dept of Electronics and Communication, R.V. College of Engineering, Bangalore, INDIA -560059

<sup>1</sup>darshanharagil.ec16@rvce.edu.in, <sup>2</sup>mahiths.ec15@rvce.edu.in

<sup>3</sup>sahanab@rvce.edu.in

**Abstract:** Kubernetes is a compact, extensible, open-source stage for overseeing containerized responsibilities and administrations, that works with both decisive setup and robotization. Kubernetes is like VMs, however having loosened up isolation properties to share the Operating System (OS) among the applications. The container conversely with VM, has its own document framework, portion of Central Processing Unit(CPU), memory, process space and much more. Kubernetes cluster is a bunch of node machines for running containerized applications. Each cluster contains a control plane and at least one node. Infrastructure Optimization is the process of analyzing and arranging the portion of cloud resources which power applications and workloads to augment the presentation and limit squander due to over-provisioning. In the paper, a "Movie Review System" web application is designed using GoLang for backend components and HTML, CSS and JS for frontend components. Using AWS, an EC2 instance is created and the web application is deployed onto EC2 and hosted in instance server. The web application is also deployed on Kubernetes locally using MiniKube tool. A performance analysis is performed for both the deployments on considering common performance metrics for both AWS EC2 / Virtual Machine (VM) and Kubernetes.

**Keywords:** Kubernetes(k8s), Kubernetes Cluster, Virtual Machines (VM), Amazon EC2, MiniKube, Web Application.

## 1. INTRODUCTION

Kubernetes is a compact, extensible, open-source stage for overseeing containerized responsibilities and administrations, that works with both decisive setup and robotization. A virtual machine, is essentially a full machine running all the components, including its own Operating System, on top of the virtualized equipment. In the year 2014 an advancement the Kubernetes which is like VMs, however having loosened up isolation properties to share the Operating System(OS) among the applications was delivered. Kubernetes initially built by Google, at present is maintained by the Cloud Native Computing Foundation(CNCF). Kubernetes furnishes with the structure to run dispersed frameworks versatily, considering the scaling and fail over of the application, giving deployment patterns and much more. Kubernetes gives the load balancing and storage orchestration, programmed binpacking, computerized rollouts and rollbacks and other network advancement related services. The Figure 1, shows the design of Kubernetes comprising of two fundamental parts: Master and Nodes. The master is divided into nodes which comprise of Pod, Container, Kubelet and Proxy. Kubernetes cluster is a bunch of node machines for running containerized applications. Each cluster contains a control plane and at least one node.

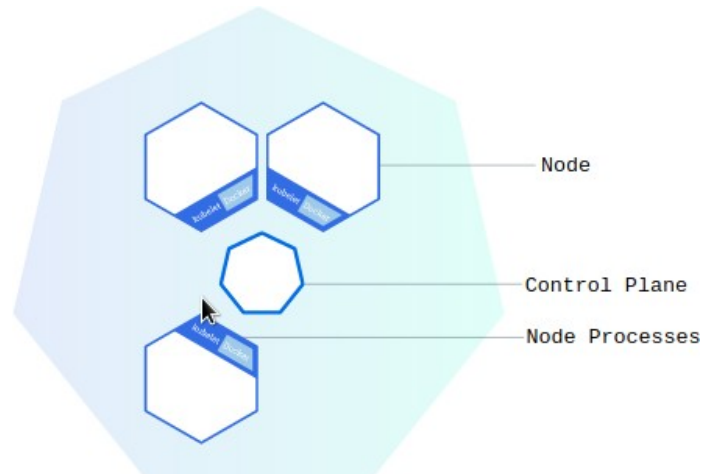


Figure 1: Kubernetes Cluster

The control plane keeps up the desired condition of cluster (like which applications are running) and the nodes run the applications and workloads. When the kubernetes is running it implies that the we are running the cluster. Each cluster comprises of individual 'Namespace' which permits to deal with multiple clusters in a similar physical cluster. Performance Optimization in the kubernetes is a significant worry recently, as kubernetes is highly extensible, open-sourceplatform for coordinating containerized workloads in cloud environment.

## 2. Literature Review

In research paper [1] the creators Nguyen Dinh Nguyen and Taehong Kim portray the Kubernetes design and the pioneer based system for keeping up predictable information stockpiling among reproductions of a stateful application in the Kubernetes group. The creators propose a pioneer political race calculation that not just works with the utilization of pioneer political decision in Kubernetes yet in addition equitably disperses the pioneers all through every one of the hubs in the group. The assessment results showed that the proposed BLD calculation can adequately adjust the quantity of pioneers among all hubs in the cluster. The viability of the BLD calculation was demonstrated through an exhibition assessment with various applications, showing that the throughput can be altogether improved by disseminating the quantity of pioneers equitably all through the hubs.

In research [2] by Ionut-Catalin Donca and group, the proposed arrangement was set off by the requirement for measure computerization, setup, execution of programming and the assembly of frameworks. Robotization of frameworks diminishes the expense of developing rehash capable science measure foundations and adds to the reproducibility of examination contemplates. The whole logical work process life cycle was rearranged by a framework introduced: programming setup, establishment and work process the board. The best technique for the inevitable arrangement of this framework is to handle the Autoscaler Cluster, which scales the cases dependent on various boundaries (CPU usage, memory use).

The creators Mulugeta Ayalew Tamiru, Johan Tordsson, Erik Elmroth and group in paper [3], express that 'Kubernetes has arisen as the accepted compartment arrangement stage in the cloud'. In the paper the creators presume that default design (CA) the autoscaler arrangements hubs at the hour of scale-out from just a single hub pool, while when arranged with hub auto-provisioning (CA-NAP) it arrangements hubs from numerous hub pools. We look at these two designs utilizing SPEC's auto-scaling execution measurements and financial expense of running the clusters, the effect of various arrangements and applications on the autoscaling execution and cost of running of a Kubernetes bunch. As the Kubernetes group autoscaler is profoundly configurable, it is fascinating to contemplate the effect of tuning extra boundaries on auto-scaling execution and cost saving.

In the paper [4] by Maria A. Rodriguez and Rajkumar Buyya, the creators have proposed the incorporated utilization of schedulers, autoscalers, and reschedulers as a component to make compartment coordination frameworks cloud-mindful. In this, the scheduler enhances the underlying arrangement of holders, the autoscaler empowers the current interest for assets to be met and underutilized or inactive hubs to be closure improve the framework's use and consequently lessen cost, lastly the rescheduler takes into account the underlying position of compartments to be updated at runtime to diminish discontinuity and combine burdens to energize better asset usage. Different rescheduling and auto-scaling system were proposed, carried out, and assessed.

Inpaper [5] by Junzo Watada, Arunava Roy and team, containerization provides many promising features like super lightweight, faster spin-up/down, efficient energy and resource utilization, impressive workload distribution capabilities, achieving server consolidation, and many more, but at the same time it has few major problems such as weaker isolation, higher chance of container sprawl, lack of capable tools for container orchestration and cross-platform supports and container portability limitations. Unikernels provide VM-like isolation with significantly small footprints along with extremely fast booting. the recent technological developments in both VMs and lightweight virtualization, the increasing industrial acceptance of containerization will soon be challenged as the need of the day has shifted to the advantages.

In research [6] by Emiliano Casalicchio, propose and assess the presentation of KHPAA an improved variant of the KHPA auto-scaling calculation. The proposed arrangement influence, in a straightforward way, total utilization estimates rather than family member. KHPAA can be connected any current framework coordinated with Kubernetes without the need to change the framework arrangement. The presentation examination shows that the utilization of total measurements permits to appropriately control the application reaction time and to keep it beneath edges presented by administration level targets.

In paper [7] by Isam Mashhour Al Jawarneh and group, have examined compartment orchestration structure contrasting highlights and administrations offered by holder orchestrators. We have likewise planned a strong arrangement of measurements to think about their exhibitions at scheduling and administration the executives layers and we accept those measurements are reusable for a wide range of holder orchestrators. At last, we picked four agent holder orchestrators, specifically, Docker Swarm, Kubernetes, and Apache Mesos.

The authors Victor Medel, Omer Rana and team in paper [8], have outlined performance model for Kubernetes based deployment. The authors use a benchmark-based approach to better characterize behavior of a Kubernetes system (using Containers). The proposed reference model can be used as support for: capacity planning and resource management is outlined and application design.

### **3. Architectures and Models**

#### **3.1 Architecture of Web Application**

Web application is a framework which outlines the key internal and external components, the cooperation of components like UIs, databases, CSV files and middleware systems. The web application architecture approaches and associates dissimilar components to improve web insight. The web application architecture comprises of segments which are categorized into two primary classes: user interface app and structural components. The Figure 2 depicts the fundamental architecture of a web application.

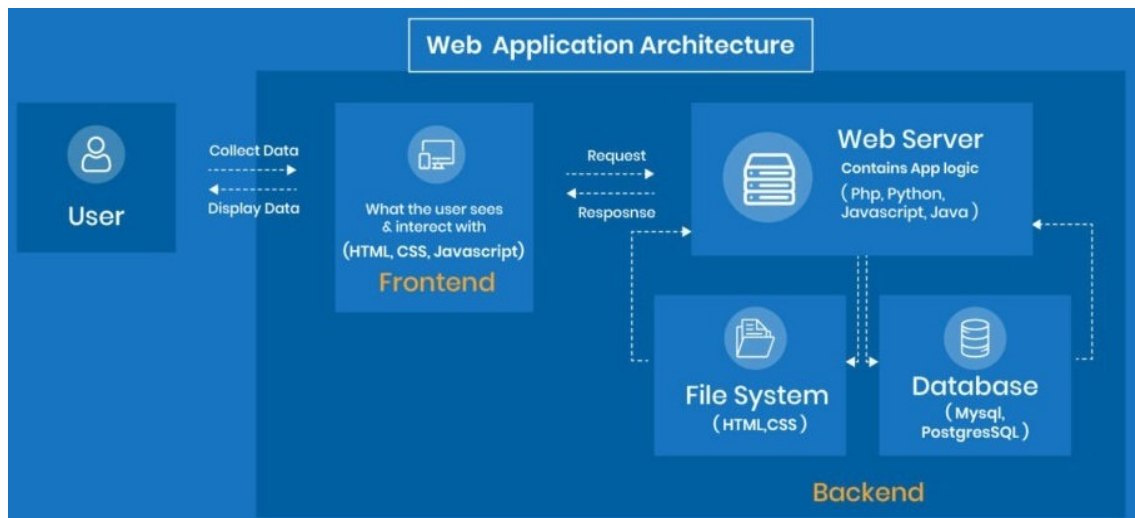


Figure 2: Basic Architecture of Web Application

### User Interface Components

User Interface Components are related to settings, configurations and display of web applications. These components are responsible for creating and improving the experience and interfaces of the web application. These contain a number of components for instance, dashboards, configuration settings, statistical data, layouts, etc.

### Structural Components

The structural components are responsible for fabricating the performance of web application. These components allow users to coordinate with the application. The structural components include:

#### The Client or Web Browser

The users can collaborate and communicate with web applications using this platform. HTML, CSS, and JavaScript are the programming languages that are valuable to shape this component.

#### CSV Files or Database

The CSV offers relevant data or business logic that is overseen and stored by web application server. The database stores, retrieves, and delivers the data.

#### Web Application Server

The web server manages the fundamental hub that upholds multi-layer applications and business logic. The server is usually built using PHP, GoLang, Python, .NET, Java and Node.js. In the project the web server is designed using the 'GoLanguage'.

## 3.2 Kubernetes Architecture

Kubernetes follows the customer worker design where we have introduced on one machine and the hub on discrete Linux machines. It uses a master to manage Docker containers across multiple Kubernetes nodes. A master and its controlled nodes constitute a Kubernetes cluster. The user can deploy an application in Docker containers via the Kubernetes master. The Figure 3 is the Kubernetes architect consisting on Master and Worker nodes.

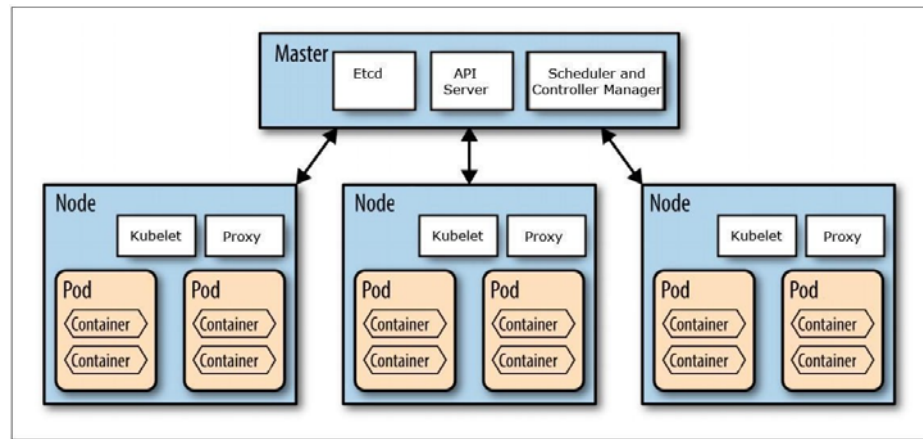


Figure 3: Kubernetes Architecture

### Kubernetes Master Components

Kubernetes master is responsible for managing the cluster, coordinating activities inside cluster and communication with worker nodes to keep the application and Kubernetes running. The Master node is the entry point of all administrative tasks. The components of Kubernetes Master:

#### API Server

The users can interact and communicate with web applications via this platform. HTML, CSS, and JavaScript are the programming languages that are valuable to build this component. Programming Interface (API) server is the section point for all the REST components used to control the cluster. All the regulatory tasks are carried out by Programming Interface (API) server within the master node. To make, erase, update or show in Kubernetes object it needs to go through this Programming Interface (API) server. Programming Interface (API) server approves and configures the Programming Interface (API) objects like ports, administrations, replication, regulators and arrangements and it is responsible for exposing Programming Interface for every operation, 'kubectl' is utilized to communicate with API tools. (API)

#### Scheduler

Scheduler is a service responsible for distributing the workload. Scheduler is liable for tracking the work load utilized at each worker node and placing the workload on which the resources are accessible. The scheduler is responsible for scheduling pods across available nodes depending on the requirements referenced in configuration document. The scheduler is liable for workload utilization and allocating pod to a new node.

#### Controller Manager

Controller Manager, also known as controllers, it is a daemon that runs in non-terminating loop and collects and sends the information to API server. The key regulators are replication regulator, endpoint regulator, namespace regulator, and administration account regulator.

#### etcd

etcd is a distributed key-value lightweight database. In Kubernetes, it is a central database for storing the current cluster state at any point of time and also used to store the configuration details such as subnets, etc. The etcd is written in GoLang.

### Kubernetes Worker Components

Kubernetes Worker contains all necessary services to manage the networking between containers, communication with the master and assign resources to scheduled containers. The components of Kubernetes Worker:

Kublet

Kublet is a primary node that communicates with the master and executes worker nodes inside the cluster. Kublet gets the pod specifications through the Application Programming Interface (API) server and executes the container associated with pods and ensure the containers are running and healthy.

Kube-Proxy

Kube Proxy is the core networking component inside the kubernetes cluster. Kube Proxy is responsible for maintaining the entire network configuration. Kube-Proxy maintains the distributed network across all the nodes, pods and the containers and also exposes the services across outside world. Kube-Proxy acts as a network proxy and load balancer for a service on a single worker node and manages the network routing for Transmission Control Protocol (TCP) and User Datagram Protocol (UDP) packets.

Pods

Pod is a gathering of containers that are deployed together on a similar host. The pods can deploy multiple dependent containers together so it acts as a wrapper around these containers so we can interact and manage the containers primarily through pods.

#### 4. Methodology

The web application is the integral part of the project. 'Movie Review System', a web application based on reviewing is designed and built. The web application consists of a Login page, Sign Up page, Home page and a Review page. A user is first directed to the Login page where he/she can directly login to their account using the user specific username and password if he/she has already created an account else the user can click on the Sign Up button on the top left corner of the Login page and create an account. On successful creation of account and signing to the account, the user gets a set of movies to choose from, for which he/she has to review and type comments and click on Submit. After submitting review, the user needs to scroll to the bottom of the page and Log Out of his account, which he/she is redirected to the Log In page. The Figure 4 is a flowchart showing the complete methodology.

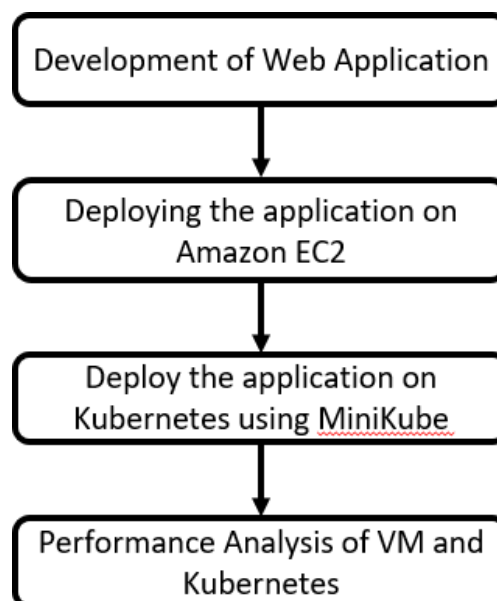


Figure 4: Complete Methodology



#### 4.1 Deployment of Web Application on Amazon EC2

To launch a web application on Amazon EC2 instance, user needs to first create an account in Amazon Web Services in order to sign into AWS Management Console. The user needs to perform following steps:

1. Choose EC2 from AWS Console and select Amazon Linux 2 Amazon Machine Image (AMI) as operating system.
2. Selecting the instance type: 't2.micro' (Free tier) from the available set of instances.
3. Assigning a Public IP address, adding Storage, Configuring Security groups which are unique to each instance.
4. On assigning all the above mentioned components, launching the instance.
5. Creating a Key pair which is unique to each instance.
5. Connecting to the instance, installing and connecting to the Apache web server.
6. After getting connected to the server, launching the web application on the AWS server.

#### 4.2 Deployment of Web Application on Kubernetes

To deploy a web application on Kubernetes locally, user needs to use MiniKube and Katacoda. Katacoda is used as it provides free and in-browser Kubernetes environment. While deploying the web application on Kubernetes the user needs to follow the steps:

1. Creating a MiniKube Cluster (locally): by opening Kubernetes dashboard in the browser in minikube. Then a Katacoda environment is selected the port as 30000.
2. Creating a Deployment: using Kubectl, creating deployment that manages the pods and checking and setting the kubectl configurations. Deployment of application checks the health of Pod and restarts the Pod's if they get terminated when being used.
3. Creating a Service: Initially the Pod is accessible only by its internal IP address, so Kubernetes service is performed to make the Pod accessible from outside the Kubernetes virtual network. The application codes only listens to client requests from 'TCP port 8080'. The Katacoda environment automatically allots a 5-digit unique port number for every Kubernetes service.
4. On completion of above mentioned tasks, the application opens on browser and its response, which is noted.

## 5. Results and Discussion

### Performance Analysis

The web application is deployed on Amazon EC2, a Cloud Compute Service provided by Amazon Web Services (AWS) and on Kubernetes using MiniKube. On hosting the application in EC2 server and locally on Kubernetes, the performance analysis is carried out using the Prometheus and Grafanna tools. The performance metrics considered are: CPU Utilization, Memory Footprints, Network Traffic, Time Synchronization and Entropy. The Figures 5.1 depicts the basic CPU utilization graphs and the Figure 5.2 shows the exact usage percentage of CPU by user, system and when at idle state.

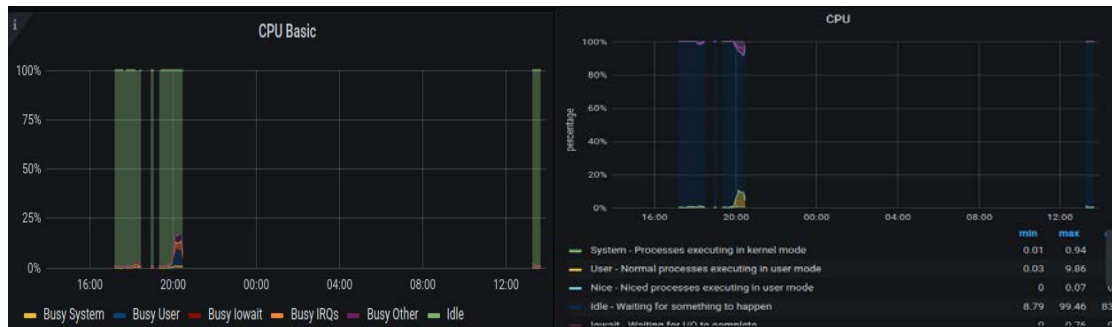


Figure 5.1: Basic CPU Utilization

Figure 5.2: Detailed CPU Usage

The Figure 5.3 shows the exact memory consumed when the user accesses the application and the Figure 5.4 shows the memory consumed when user switched between pages.



Figure 5.3: Memory Usage Status

Figure 5.4: Memory used for Pages Accessed

The Figures 5.5 and 5.6 shows the time synchronized/elapsed between one LogIn and LogOut by the user to his/her account in the web application.

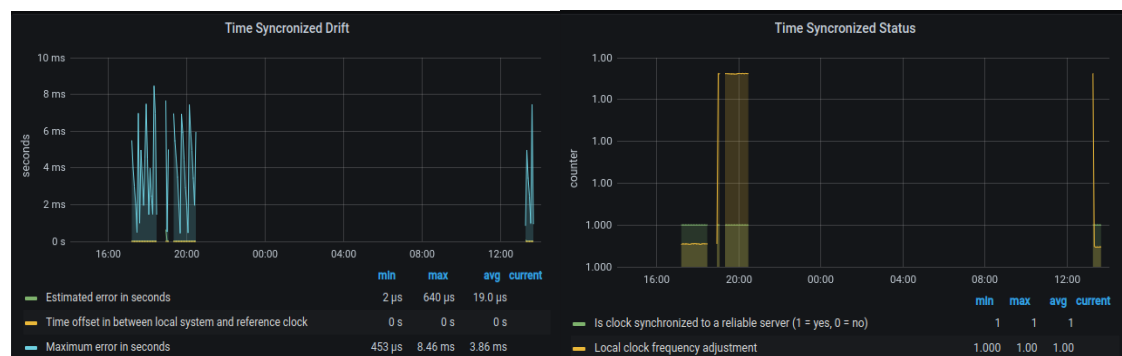


Figure 5.5: Time Synchronization

Figure 5.6: Time Usage Status



The Figure 5.7 depicts the Entropy of the web application i.e, the number of remote users accessing the web application at a given time.

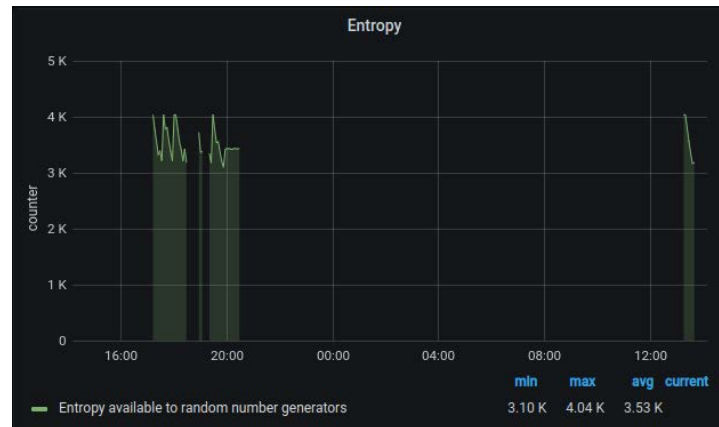


Figure 5.7: Entropy of web application

The Figure 5.8 shows Network traffic i.e, the amount of data which is moving across the web application/computer network. The Network traffic is analyzed in terms of data packets.

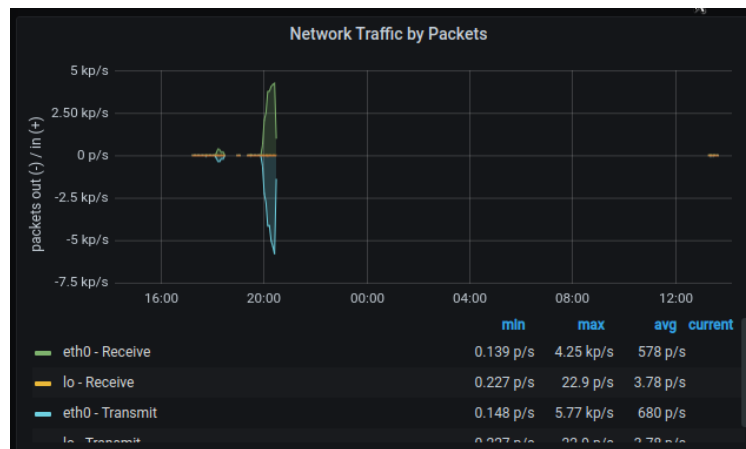


Figure 5.8: Network Traffic in terms of Data Packets

## 6. Conclusion

A 'Movie Review System', a web application based reviewing is designed and built. The application uses GoLang for backend and HTML, CSS and JS for frontend. The web application consists of a Login page, Sign Up page, Home page and a Review page. A user has to Log In to his/her account, select a movie, review and submit and then Log Out of his/her account.

Amazon Web Services, an open-source platform for a variety of cloud computing operations is used to host the web application on the server. The web application is deployed on the Amazon EC2 instance (Virtual Machine) and the instance is hosted on Amazon server. Using Minikube, the web application is deployed on Kubernetes locally. The performance metrics considered in common are noted for both the Virtual Machine and Kubernetes deployments.

The performance analysis is performed to conclude which among the Virtual Machine and Kubernetes performs better. The performance analysis is carried out considering the

performance metrics in common to both the VM and Kubernetes. Based on the results obtained after performance analysis, we conclude the Kubernetes performs better in terms of CPU utilization, Network traffic and Memory usage. Hence, from this study we conclude that the Kubernetes is better when compared to regular Virtual Machine when it comes to Load Balancing the Cloud Environments.

## 7. References

1. N.D.Nguyen and T.Kim, "Balanced leader distribution algorithm in Kubernetes clusters," *Sensors*, vol. 21, no. 3, 2021, ISSN: 1424-8220. DOI: 10.3390/s21030869. [Online]. Available: <https://www.mdpi.com/1424-8220/21/3/869>. Sahoo, Doyen Liu, Chenghao Hoi, Steven. (2017). Malicious URL Detection using Machine Learning: A Survey.
2. C.Donca, C.Corches, O.Stan, and L.Miclea, "Autoscaled rabbitmq Kubernetes clusters on single-board computers," in *2020 IEEE International Conference on Automation, Quality and Testing, Robotics (AQTR)*, 2020, pp. 1–6. DOI: 10.1109/AQTR49680.2020.9129886.
3. J.M.A.Tamiru, J.Tordsson, E.Elmroth, and G.Pierre, "An experimental evaluation of the Kubernetes cluster autoscaler in the cloud," in *2020 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, 2020, pp. 17–24. DOI: 10.1109/CloudCom49646.2020.00002. V.M. Patro, M.R. Patra, Augmenting Weighted Average with Confusion matrix to Enhance classification accuracy. *Trans. Mach., Learn. Artif Intell.* 2(4), 77–91 (2019)
4. M. A. Rodriguez and R. Buyya, "Containers orchestration with cost-efficient autoscaling in cloud computing environments," *ArXiv*, vol. abs/1812.00300, 2018.
5. J.Watada, A.Roy, R.Kadikar, H.Pham, and B.Xu, "Emerging trends, techniques and open issues of containerization: A review," *IEEE Access*, vol. 7, pp. 152443–152472, 2019. DOI: 10.1109/ACCESS.2019.2945930.
6. E.Casalicchio, "A study on performance measures for auto-scaling CPU-intensive containerized applications," *Cluster Computing*, vol. 22, Sep. 2019. DOI: 10.1007/s10586-018-02890-1.
7. A.Pereira Ferreira and R.Sinnott, "A performance evaluation of containers running on managed Kubernetes services," in *2019 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, 2019, pp. 199–208. DOI: 10.1109/CloudCom.2019.00038.
8. I.M.A.Jawarneh, P.Bellavista, F.Bosi, L.Foschini, G.Martuscelli, R.Montanari, and A.Palopoli, "Container orchestration engines: A thorough functional and performance comparison," in *ICC 2019-2019 IEEE International Conference on Communications (ICC)*, 2019, pp. 1–6. DOI: 10.1109/ICC.2019.8762053.